

Citation for published version:

Webster, J, Exley, J, Copestake, J, Davies, R & Hargreaves, J 2018, 'Timely evaluation in international development', *Journal of Development Effectiveness*, vol. 10, no. 4, pp. 482-508.
<https://doi.org/10.1080/19439342.2018.1543345>

DOI:

[10.1080/19439342.2018.1543345](https://doi.org/10.1080/19439342.2018.1543345)

Publication date:

2018

Document Version

Peer reviewed version

[Link to publication](#)

Publisher Rights

Unspecified

This is an Accepted Manuscript of an article published by Taylor & Francis in *Journal of Development Effectiveness* on 27/11/2018, available online: <http://www.tandfonline.com/10.1080/19439342.2018.1543345>

University of Bath

Alternative formats

If you require this document in an alternative format, please contact:
openaccess@bath.ac.uk

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

PROOF COVER SHEET

Journal acronym: RJDE

Author(s): Jayne Webster, Josephine Exley, James Lewis, James Copestake,
Rick Davies and James Hargreaves

Article title: Timely evaluation in international development

Article no: 1543345

Enclosures: 1) Query sheet
2) Article proofs

Dear Author,

1. Please check these proofs carefully. It is the responsibility of the corresponding author to check these and approve or amend them. A second proof is not normally provided. Taylor & Francis cannot be held responsible for uncorrected errors, even if introduced during the production process. Once your corrections have been added to the article, it will be considered ready for publication.

Please limit changes at this stage to the correction of errors. You should not make trivial changes, improve prose style, add new material, or delete existing material at this stage. You may be charged if your corrections are excessive (we would not expect corrections to exceed 30 changes).

For detailed guidance on how to check your proofs, please paste this address into a new browser window: <http://journalauthors.tandf.co.uk/production/checkingproofs.asp>

Your PDF proof file has been enabled so that you can comment on the proof directly using Adobe Acrobat. If you wish to do this, please save the file to your hard disk first. For further information on marking corrections using Acrobat, please paste this address into a new browser window: <http://journalauthors.tandf.co.uk/production/acrobat.asp>

2. Please review the table of contributors below and confirm that the first and last names are structured correctly and that the authors are listed in the correct order of contribution. This check is to ensure that your name will appear correctly online and when the article is indexed.

Sequence	Prefix	Given name(s)	Surname	Suffix
1		Jayne	Webster	
2		Josephine	Exley	
3		James	Lewis	
4		James	Copestake	
5		Rick	Davies	
6		James	Hargreaves	

Queries are marked in the margins of the proofs, and you can also click the hyperlinks below.

Content changes made during copy-editing are shown as tracked changes. Inserted text is in **red font** and revisions have a red indicator ▲. Changes can also be viewed using the list comments function. To correct the proofs, you should insert or delete text following the instructions below, but **do not add comments to the existing tracked changes**.

AUTHOR QUERIES

General points:

1. **Permissions:** You have warranted that you have secured the necessary written permission from the appropriate copyright owner for the reproduction of any text, illustration, or other material in your article. Please see <http://journalauthors.tandf.co.uk/permissions/usingThirdPartyMaterial.asp>.
2. **Third-party content:** If there is third-party content in your article, please check that the rightsholder details for re-use are shown correctly.
3. **Affiliation:** The corresponding author is responsible for ensuring that address and email details are correct for all the co-authors. Affiliations given in the article should be the affiliation at the time the research was conducted. Please see <http://journalauthors.tandf.co.uk/preparation/writing.asp>.
4. **Funding:** Was your research for this article funded by a funding agency? If so, please insert 'This work was supported by <insert the name of the funding agency in full>', followed by the grant number in square brackets '[grant number xxxx]'.
5. **Supplemental data and underlying research materials:** Do you wish to include the location of the underlying research materials (e.g. data, samples or models) for your article? If so, please insert this sentence before the reference section: 'The underlying research materials for this article can be accessed at <full link> / description of location [author to complete]'. If your article includes supplemental data, the link will also be provided in this paragraph. See <http://journalauthors.tandf.co.uk/preparation/multimedia.asp> for further explanation of supplemental data and underlying research materials.
6. The **CrossRef database** (www.crossref.org/) has been used to validate the references. Changes resulting from mismatches are tracked in **red font**.

AQ1 Please provide a short biography.

AQ2 Please provide missing City/Country for the affiliation(a-d).

AQ3 Please note that the Funding section has been created from information provided through CATS. Please correct if this is inaccurate.

AQ4 The disclosure statement has been inserted. Please correct if this is inaccurate.

AQ5 The CrossRef database (www.crossref.org/) has been used to validate the references. Mismatches between the original manuscript and CrossRef are tracked in red font. Please provide a revision if the change is incorrect. Do not comment on correct changes.

AQ6 Please provide missing Publisher location for the "Barder and Ramalingam, 2012" references list entry.

AQ7 Please provide missing Publisher location for the "Barr, 2015" references list entry.

- AQ8** Please provide missing Publisher location for the "Beach and Pedersen, 2013" references list entry.
- AQ9** Please provide missing Publisher location/Publisher name for the "Beebe, 2001" references list entry.
- AQ10** Please provide missing Publisher location/Publisher name for the "Butler, 1995" references list entry.
- AQ11** Please provide missing Publisher location for the "Davies, 1996" references list entry.
- AQ12** Please provide missing Publisher location for the "Davies, 2014" references list entry.
- AQ13** Please provide missing Publisher location \Name for the "Davies and Dart, 2005" references list entry.
- AQ14** Please provide missing Publisher location for the "Davies et al., 2016" references list entry.
- AQ15** Please provide missing Publisher location for the "DFID, 2012" references list entry.
- AQ16** Please provide missing Publisher location for the "Eirich and Morrison, n.d." references list entry.
- AQ17** Please provide missing Publisher location for the "Fereday, 2015" references list entry.
- AQ18** Please provide missing Publisher location for the "Goertz and Mahoney, 2012" references list entry.
- AQ19** Please provide missing Publisher location for the "Green, 2015" references list entry.
- AQ20** Please provide missing Publisher location for the "Global Pulse, 2012" references list entry.
- AQ21** Please provide missing Publisher location for the "Hubbard, 2010" references list entry.
- AQ22** Please provide missing Publisher location for the "IPA, 2016" references list entry.
- AQ23** Please provide missing Publisher location for the "Jones and Hearn, 2009" references list entry.
- AQ24** Please provide missing Publisher location for the "Karlan, 2017" references list entry.
- AQ25** Please provide missing Publisher location for the "Ladner, 2015" references list entry.
- AQ26** Please provide missing Publisher location/Publisher name for the "ODI, 2009" references list entry.
- AQ27** Please provide missing Publisher location for the "Patton, 2008" references list entry.
- AQ28** Please provide missing Publisher location for the "Pawson, 2013" references list entry.
- AQ29** Please provide missing Publisher location/Publisher name for the "Pawson and Tilley, 2004" references list entry.
- AQ30** Please provide missing Publisher location for the "Rio et al., 2015" references list entry.
- AQ31** Please provide missing page number for the "Schünemann and , 2015" references list entry.

- AQ32** Please provide missing volume number for the “Smutylo, 2005” references list entry.
- AQ33** Please provide missing Publisher location for the “Valters et al., 2016” references list entry.
- AQ34** Please provide missing Publisher location for the “Wilson-Grau and Britt, 2012” references list entry.

PROOF ONLY

How to make corrections to your proofs using Adobe Acrobat/Reader

Taylor & Francis offers you a choice of options to help you make corrections to your proofs. Your PDF proof file has been enabled so that you can mark up the proof directly using Adobe Acrobat/Reader. This is the simplest and best way for you to ensure that your corrections will be incorporated. If you wish to do this, please follow these instructions:

1. Save the file to your hard disk.
2. Check which version of Adobe Acrobat/Reader you have on your computer. You can do this by clicking on the "Help" tab, and then "About".

If Adobe Reader is not installed, you can get the latest version free from <http://get.adobe.com/reader/>.

3. If you have Adobe Acrobat/Reader 10 or a later version, click on the "Comment" link at the right-hand side to view the Comments pane.

4. You can then select any text and mark it up for deletion or replacement, or insert new text as needed. Please note that these will clearly be displayed in the Comments pane and secondary annotation is not needed to draw attention to your corrections. If you need to include new sections of text, it is also possible to add a comment to the proofs. To do this, use the Sticky Note tool in the task bar. Please also see our FAQs here: <http://journalauthors.tandf.co.uk/production/index.asp>.

5. Make sure that you save the file when you close the document before uploading it to CATS using the "Upload File" button on the online correction form. If you have more than one file, please zip them together and then upload the zip file.

If you prefer, you can make your corrections using the CATS online correction form.

Troubleshooting

Acrobat help: <http://helpx.adobe.com/acrobat.html>

Reader help: <http://helpx.adobe.com/reader.html>

Please note that full user guides for earlier versions of these programs are available from the Adobe Help pages by clicking on the link "Previous versions" under the "Help and tutorials" heading from the relevant link above. Commenting functionality is available from Adobe Reader 8.0 onwards and from Adobe Acrobat 7.0 onwards.

Firefox users: Firefox's inbuilt PDF Viewer is set to the default; please see the following for instructions on how to use this and download the PDF to your hard drive:

http://support.mozilla.org/en-US/kb/view-pdf-files-firefox-without-downloading-them#w_using-a-pdf-reader-plugin

ARTICLE



Timely evaluation in international development

Jayne Webster^a, Josephine Exley^b, James Lewis^c, James Copestake^d, Rick Davies^e
and James Hargreaves^b

AQ1

AQ2

^aCentre for Evaluation and Disease Control Department, London School of Hygiene and Tropical Medicine (LSHTM); ^bCentre for Evaluation and Department of Social and Environmental Health Research, LSHTM; ^cCentre for Evaluation and Department of Infectious Disease Epidemiology, LSHTM; ^dCentre for Development Impact, University of Bath; ^eIndependent Consultant, UK

ABSTRACT

Impact and process evaluations are increasingly used in international development; however they are generally retrospective in outlook. A more timely approach to evaluation aims to identify necessary, feasible and effective changes during a programme or intervention's lifetime. This paper aims to identify, categorise, describe and critically appraise methods to support more timely evaluation in international development. Potential methods were identified through scoping seminar, public symposium, targeted review of the literature, and the authors' own experiences and opinions. Findings from the different data sources were reviewed collectively by the author group and triangulated to develop an analytical framework. We identified four purposes of timely evaluation for international development, and critiqued the use of these approaches against four dimensions of timeliness and flexibility. Whilst we found significant interest in more timely approaches to evaluation in international development, there was a dearth of published empirical evidence upon which to base strong recommendations. There is significant potential for timely evaluation to improve international development outcomes. New approaches to mixing and adapting existing methods, together with new technologies offer increased potential. Research is needed to provide an empirical evidence base upon which to further develop the application, across sectors and contexts, of timely evaluation in international development.

ARTICLE HISTORY

Received 26 July 2018
Accepted 30 October 2018

KEYWORDS

Outcome evaluation; impact evaluation; adaptive learning; programme improvement

Introduction

Outcome evaluations assess the impact of a specified set of actions, constituting a programme or intervention, on its intended outcomes. Such evaluations ask: what effect did this action have on these outcomes (often in comparison with some other action). Process evaluations seek to explain how and why such impacts, did or did not, come about (Moore et al. 2015). They assess how implementation of a programme happened, whether hypothesised causal pathways were activated and identify contextual factors that acted as barriers or facilitators to either implementation, effectiveness, or both. Such evaluations are essential for informing future policy decisions, but many of the questions typically addressed are, by their nature, retrospective in outlook.

Dealing with the uncertainty and complexity inherent in international development settings requires a flexible approach to the design and implementation of programmes. Flexibility is needed across time (for example, changing activities or shifting priorities over time) and space (for

example, adapting an approach to different settings and contexts), and happens at multiple speeds (for example, daily fine tuning of specific activities, annual changes in budget allocations and longer-term priority setting) (Barder and Ramalingam 2012; Gamble 2006; Ladner 2015; Valters, Cummings, and Nixon 2016; Walji and Vein 2013). The *Doing Development Differently* manifesto highlights that, among other things, to be successful development programmes need to 'merge design and implementation' by undertaking 'rapid cycles of planning, action, reflection and revision' and 'manage risk by making small bets; pursuing activities with promise and dropping others' (DDD 2014).

Evaluations have a role to play in supporting the Doing Development Differently agenda by generating evidence to inform action during a programme's life cycle; from design to the selection, refinement and testing of interventions. Where knowledge is high about what is likely to work, evaluation can test whether the intervention is having the anticipated effect and support, and test modifications over time. Where it is less clear what intervention might work, interventions need to be developed and options tested either sequentially or in parallel (Green 2015; O'Donnell 2016).

Despite there being a number of existing approaches and methods to incorporating evidence based decision making into programmes, there has been scant focus on, or critique of, 'timeliness' and the suitability of evaluation methods within flexible or adaptive international development programmes. We aim to review and critically appraise evaluation methods to support a more 'timely' approach to evaluations of international development programmes. To support this critical appraisal we define a 'timely' approach to evaluation and consider purposes of the evaluation and dimensions of the methods required for timely application and decision making. To guide evaluators we propose a framework to support the selection of methods, or mixes of methods, needed to address particular evaluation questions at different stages of a programme's cycle.

Methods

Our review and critique of methods for timely evaluation included: a scoping seminar and public symposium to identify methods from the perspectives of academics, programme designers and programme evaluators; a review of approaches and methods used to evaluate international development programmes; and a critique of methods against a timely evaluation framework.

Scoping seminar and public symposium

The scoping seminar on 'real time evaluations for programme improvement' took place in June 2017 at the London School of Hygiene and Tropical Medicine (LSHTM) to harness the ideas and experiences of members of LSHTM's Centre for Evaluation. The seminar was attended by approximately 30 members from a range of disciplines within public health. The seminar included six speed talks and a group discussion. The public symposium held in November 2017 was attended by 142 people and included three sessions on: doing, evaluating and critiquing timely evaluations for programme improvement. Presentations were given by eight speakers. During the event we engaged with participants through breakout sessions and technology. We had an active twitter discussion (#timelyeval) and used slido.com for participants to submit questions/comments during presentations. Both events were recorded and in drafting this manuscript we listened back to the recordings and took notes. Through the presentations and group discussions at the two events we collated a list of potential methods to examine in more detail.

Review of approaches and methods to evaluate international development programmes

The literature review consisted of two components. First, following the scoping event, we undertook a targeted review using a snowballing technique to identify specific methods that have been used in evaluations of adaptive learning approaches in development settings (Wohlin 2014). Based

Table 1. Search terms.

Search	Terms (title/abstract/key word)
1	'Adaptive learn*' OR 'continuous evaluat*' OR 'developmental evaluat*' OR 'experiential learn*' OR 'feedback' OR 'formative evaluat*' OR 'real time evaluat*' OR 'Problem Driven Iterative Adaptation'
2	Humanitarian OR International Development
3	1 AND 2

on the scoping seminar, we developed a set of search terms (Table 1). Searches were run in PubMed and Web of Science. The reference list of relevant literature was screened, and we undertook forward citation searching in Google scholar. Second for the specific methods identified during the two events, targeted searches were run in google, google scholar, PubMed and Web of Science to identify examples of where the methods had been used in international development contexts.

Critique of methods against the timely evaluation framework

We developed a framework for timely evaluation of international development programmes and interventions based on our interpretation of the discussion at the scoping event, public symposium and review of the literature. We critiqued examples of the methods against the timely evaluation framework.

Results

Based on the discussions at the scoping event and public symposium, we defined a timely approach to evaluation as 'the use of evaluation methods before or during the course of an international development programme or intervention to provide evidence for decision making on design, adaptation or refinement at a time when these changes can plausibly lead to the improvements needed, and when implementers and stakeholders can effectively carryout and benefit from the changes'. This definition highlights the interconnected nature of timeliness and flexibility, which we expand on below.

During the internal and external events, participants highlighted an array of existing approaches that they considered encapsulated aspects of a timely approach to evaluation, including programme cycles, quality improvement, rapid cycle evaluations and developmental evaluations. Additional related approaches were identified through the literature review. At their core these approaches aim to generate more timely evidence over a programme or interventions life cycle and respond to changing and evolving priorities. The complete list of approaches identified is listed in Table 2.

The approaches listed in Table 2 often consist of a number of different methods. The challenge for evaluators is to identify suitable methods that can be used over varying timeframes to answer different evaluation questions at different time points as the programme unfolds. We summarise

Table 2. Approaches for timely evaluation and adaptive learning.

Accountable aid, action research, active research, adaptive development, adaptive learning, adaptive management, adaptive programming, adaptive strategy, agile working practices, appreciative inquiry, augmented feedback, behaviour centred design/human centred design, better programme delivery, citizen engagement, collaborating learning and adapting, complexity thinking, constituent voice, continuous evaluation, continuous improvement, creative design process, developmental evaluation, dynamic adaptive pathways, experiential learning, extrinsic feedback, feedback loops, feedback mechanisms, formative evaluation, iterative inquiry framework, iterative evaluation process, knowledge of results feedback, lean startup learning culture/system, model for improvement, nimble evaluations, performance management, plan-do-study-act cycle, problem driven iterative adaptation, problem based iterative adaptation, quality improvement, rapid assessment/rapid assessment process/rapid assessment methodology, rapid-cycle assessment, rapid cycle evaluation, rapid cycle quality improvement, rapid evaluation (and assessment) methods, rapid feedback evaluation, rapid qualitative enquiry, real time adaption, real time evaluation, social learning, strategy testing, utilisation focused evaluation

the methods identified through the scoping seminar, symposium and literature review in Table 3. The methods are both quantitative and qualitative, retrospective and prospective in their outlook, involve differing levels of technical skills in their analysis, and are generally applied at different stages of and time points within programmes for different purposes.

Framework for timely approach to evaluation

120

To support the selection of methods, we conceptualise a timely approach to evaluation around an analytical framework (Figure 1). The framework consists of four overarching purposes and four timeliness and flexibility dimensions. The framework recognises that methods can be used at different time points in the programme cycle and that the methods have different levels of flexibility that will make them more or less suitable in specific settings and contexts.

125

Purpose

The overarching purposes identified are: support design, identify problems, test potential solutions and explain the outcomes.

Support design. Of an intervention or package of interventions within a programme conducted prior to and/or during implementation. Where data are collected prior to implementation the purpose is to make suggestions about what interventions should be implemented and how; or to determine modifications needed to a pre-existing intervention to implement in a new context. Where a programme or intervention is already running the purpose is to explore why an anticipated change might not have occurred and identify new interventions, changes to intervention design, or implementation strategies for existing interventions in reaction to identified problems.

130

135

Identify problems. Where an intervention or programme is running the purpose is to monitor the status of implementation and identify problems that might need to be responded to. Monitoring may include all or a selection of components of a programme. Achievements are assessed against expectations that may be defined pre- or during implementation.

Test potential solutions. Where need has been identified, the purpose is to test potential options and explain why they do or do not succeed in achieving the changes required. That is, evaluating whether particular interventions or course corrections are successful in meeting their stated objectives, or are comparatively better than other options, at a given time point during the programme.

140

Explain the outcomes. Where problems in implementation or achievements have been identified and options/solutions are tested, it is important to understand and explain the outcomes. Understanding how the tested solutions change the interventions, programmes or their implementation to facilitate improvement and increase the potential for learning.

145

The four purposes are not anticipated to proceed in a cyclical manner. For example, where a new design is identified or modification made the next step may be to test potential solutions or where a problem is identified then further research may seek to support the design of potential solutions to the problem.

150

Timeliness and flexibility dimensions

We identify four timeliness and flexibility dimensions that can be used to select between methods for specific purposes: design, speed, capacity and space. The choice of method will depend on the required level of flexibility and potential time constraints. The dimensions should be considered together as they are overlapping and exert mutual influences one to the other.

155

Table 3. Evaluation methods reviewed.

Method	Description	Use and timing	Strengths/Weaknesses/Considerations
A/B tests (also known as Nimble RCT, split tests, rapid-fire tests, bucket testing, randomised field experiments) (Dibner-Dunlap and Rathore 2016; IPA 2016; Karlan 2017; Optipedia n.d.)	Clinical study design; participants are randomly assigned to receive a variation of the same intervention. Compares the effect of the adaptations on short-term outcomes.	<ul style="list-style-type: none"> • Simultaneous testing of low-cost modifications to a programme's design or message, where changes are anticipated to result in immediate change. • Particularly useful at the design or pilot stage of a programme and for answering questions about the early stages of a programme's theory of change. • Focus on short-term outcomes and use of pre-existing data enables rapid testing of elements of a programme within a relatively short time frame. 	<ul style="list-style-type: none"> • Focus on shorter-term outcomes, such as uptake and use but does not provide insight on whether the changes had an impact on longer-term changes. • Small effect sizes as examining incremental change; requires large samples. • Relies on good quality routine/administrative data being available.
Adaptive randomised control trial (Bhatt and Mehta 2016; Kairalla et al. 2012; Korn and Freidlin 2017; Lang 2011; Villar, Bowden, and Wason 2017; Cellamare et al. 2017; Choko et al. 2017; Bothwell et al. 2018; Mahajan and Gupta 2010; Thorlund et al. 2018)	Clinical study design; compares outcomes between control and intervention group. Outcomes are analysed at predefined interim time points and modifications to the study can be implemented based on the findings of the interim analysis. Modifications are made based on pre-specified decision rules.	<ul style="list-style-type: none"> • Where not clear which interventions are most likely to be effective to achieve similar outcomes, as allow simultaneous testing of multiple experimental arms. • Ongoing learning based on interim analysis: stop or start treatment arms; adjust the study population and sample size; skew treatment allocation to those treatments that appear to be doing better. • Provides confirmatory learning at end of trial. • Reduces time by combining trial phases into a single study. 	<ul style="list-style-type: none"> • Ability to make adjustments to the intervention or trial design as data is being collected, without undermining the validity or integrity of the study. • Outcomes to be measured specified at trial outset. • Decisions made during trial based on interim-findings. • More resource intensive; requires interim data collection and more rounds of analysis than a classic RCT. • Increased trial complexity; requires sophisticated statistical techniques for the analysis. • Introducing new trial arms reduces statistical efficiency. • Potential for bias from temporal trends, for example participants recruited at early stages differ to those recruited at latter stages. • Requires population level data; routine or programmatic survey data • Requires a hypothesised casual pathway; assumes achieving one step is a necessary condition to achieving the next, for example the Theory of Change is the only route through which change can occur. • Requires an understanding of whether the population in one stage is the same as the population in the next, to ascertain whether it is a 'necessary' condition or whether other steps, not captured in the Theory of Change, might be sufficient to achieve the desired change. • Casual pathways can be modified overtime. • Does not assess causality.
Bottleneck analysis/Cascade analysis/Community or systems effectiveness/Funnel of attrition (Davies 2014; Dellicour et al. 2016; Garnett et al. 2016; O'Connell and Sharkey 2013; Rio et al. 2015; Tanahashi 1978; Webster et al. 2013; White 2013)	Quantitative analysis. Identifies the steps that link the intended beneficiaries from the actual beneficiaries. Each step is conditional on the previous one having been met and only the population left at the end of all the steps have achieved the desired outcome. The relative size of the population lost at each step might indicate where the most urgent action is needed. Analysis can be stratified to understand differences between sub-groups.	<ul style="list-style-type: none"> • Identifies component(s) of a system that limits its overall performance or capacity. • Undertaken once an intervention is running and anticipate that an impact should have occurred. • Often undertaken at a single point in time providing a snap shot of need; where routine or programme data is available analysis could be repeated to assess whether the bottlenecks identified, and size, change overtime. 	

(Continued)

Table 3. (Continued).

Method	Description	Use and timing	Strengths/Weaknesses/Considerations
Contribution analysis (Befani and Mayne 2014; Eirich and Morrison n.d.; Mayne 2008)	A structured approach to explore and estimate the relative contribution of an intervention to an outcome. Maps out ongoing activities that are being undertaken that are expected to contribute to a particular outcome. Collects diverse evidence to populate 'performance stories' against a pre-specified theory of change.	<ul style="list-style-type: none"> Used to confirm or revise a theory of change. Provide feedback on what is driving change and relative contribution of a particular intervention. Particularly useful in situations where an experimental method is not feasible. Best suited to large scale programmes 	<ul style="list-style-type: none"> Retrospective approach, little or no scope for varying how the programme is implemented. Considers the relative impact of other activities on a desired outcome.
Ecological momentary assessment. Ambulatory Assessment/ Experience Sampling (Burke et al. 2017; Shiffman, Stone, and Hufford 2008)	Longitudinal design; method for collecting data in real-time, in real world settings. Participants complete short assessments on their current experiences/ behaviours/moods/environment at multiple random moments over time. Two approaches: (1) signal-contingent recording – assessed a fixed number of times per day/week, and so forth. On a random schedule; (2) event-contingent recording – assessed following exposure to specific events.	<ul style="list-style-type: none"> Used to study psychological, behavioural, and physiological processes in the natural environment. When using mobile technology data generated in real-time. 	<ul style="list-style-type: none"> Minimises recall bias; combines actual exposure measurements with momentary-measured outcomes. Repeat sampling of same individuals allows for within- and between-participant analysis. Can examine causality between exposures and behaviours. Challenges; logistic, analytic, and interpretation problems Increasing availability of mobile technology offers increased utility
Interrupted time series analysis (Biglan, Ary, and Wagenaar 2000; Kontopantelis et al. 2015; Lopez Bernal, Cummins, and Gasparini 2018; Lopez Bernal, Cummins, and Gasparini 2017)	A quasi-experimental method (others include difference-in-difference, synthetic controls, matching, regression discontinuity); model trend in outcome before and after intervention is introduced against what would have happened is the intervention was not introduced. Any change in the level of the outcome or in the rate of change over time, compared to the model, can be interpreted as the effect of the intervention.	<ul style="list-style-type: none"> To determine the effect of an intervention implemented at a specific time point in the absence of a parallel control. More complex designs can be used in situations where intervention is stopped/reversed or with multi-component interventions where different steps are implemented at different time points. 	<ul style="list-style-type: none"> Requires a large amount of data to be collected before and after intervention is introduced at equally-spaced time intervals. Outcomes to be measured need to be pre-specified. Population under study act as own control; although analysis can also include a control group, for example from a different area. Requires programmatic or routine data.

(Continued)

Table 3. (Continued).

Method	Description	Use and timing	Strengths/Weaknesses/Considerations
Modified stepped wedge trials (Wechsberg et al. 2017)	Clinical study design; compares outcomes between control and intervention arms within each step. A modified design incorporates a period of reflection at the end of each step for example undertaking surveys/IDIs/FGDs to understand how the intervention is working, and so forth. Modifications to intervention can be implemented before the next step.	<ul style="list-style-type: none"> • Prospective; to test and adapt implementation strategies. • Ongoing learning; make sequential changes to the intervention. • Confirmatory learning at the end of trial possible to compare the effect of the overall package of interventions on the pre-specified outcomes as in the original study. 	<ul style="list-style-type: none"> • Potential bias from temporal trends, for example if participants recruited early in the trial differ to those recruited later. Adaptations made during trial based on interim analysis. • Outcomes to be measured specified at trial outset; although additional unanticipated outcomes can be explored in the 'period of reflection'. • More resource intensive; requires additional data collection between steps; time needed to undertake data collection and analysis can increase length of trial. • Can increase trial complexity • Limited evidence of use from literature. • Does not get at causality • Measures intermediate outcomes and programme impact. • Can capture unexpected outcomes as do not have to hypothesise causal pathways between activities and outcomes. • Stories collected at a single point in time so does not account for changes due to temporality. • Human resource intensive; collect stories, convene panels and feedback findings. • Suitable when inputs, activities and outputs and the causal mechanisms between them are not fully understood as does not measure pre-determined outcomes. • Can identify unintended outcomes • Tailored to project and context; findings not generalisable. • Only outcomes informant aware of captured. • Resources intensive • Participation of those who influence outcome
Most significant change (Connors et al. 2017; Davies and Dart 2005; Ho et al. 2015; Limato et al. 2018; White and Phillips 2012)	Participatory qualitative method; use qualitative methods to collect programme beneficiaries' stories of recent significant change in their lives and the key activities they think led to these changes. Panel of stakeholders select what they consider to be the most significant stories, to arrive at a reduced set of changes.	<ul style="list-style-type: none"> • Retrospective; undertaken when anticipate some impact should have occurred. • Can be undertaken on an ongoing basis throughout the project cycle to reveal changes in stakeholder's perspectives at different time points. • Useful in contexts where programme already running or highly complex setting and not clear what impact may have. 	
Outcome harvesting (Wilson-Grau 2015; Wilson-Grau and Britt 2012)	Participatory approach; stakeholders collect evidence of what has changed, then work backwards to determine whether and how an intervention contributed to these changes. Draws on IDIs and surveys.	<ul style="list-style-type: none"> • Provides retrospective learning about what was achieved and how, regardless of whether it was pre-defined. • Requires an understanding of when might anticipate change to have occurred. • Useful in context where relationship between cause and effect are not fully understood. 	

(Continued)

Table 3. (Continued).

Method	Description	Use and timing	Strengths/Weaknesses/Considerations
Outcome mapping (Earl, Carden, and Smutylo 2001; Jones and Hearn 2009; ODI 2009; Research to Action 2012; Smutylo 2005)	Focuses on changes in behaviour, relationships, actions and activities of the people, groups and organisations it works with directly ('boundary partners') and how far changes contributed to outcomes. Consists of three stages: (1) intentional design – to establish consensus on intended changes; (2) outcome and performance monitoring – uses journals to chart changes in the indicators defined; (3) evaluation planning – helps the programme identify evaluation priorities and develop an evaluation plan.	<ul style="list-style-type: none"> Used at the project outset to identify activities and the individuals, groups, organisations need to work with to realise intended outcomes. 	<ul style="list-style-type: none"> Process is more intensive because it requires meaningful participation from boundary partners. Findings will be context-specific. Participatory approach means individuals involved in the project gain an understanding of their role in ensuring programme is a success. Challenges in participatory approaches of unequal power relationships.
Positive deviants (Andrews 2015; Busza et al. 2017; Positive Deviance Initiative 2017)	Explores an individual's or group's, behaviours or characteristics that have enabled them to succeed when the majority of peers have failed when faced with similar challenges, constraints, and so forth. These cases can be identified by both participatory means and more quantitative modelling approaches.	<ul style="list-style-type: none"> To discover the inputs and activities that have driven success and thus identify solutions that can be tested elsewhere. 	<ul style="list-style-type: none"> Small sample size Reflects perspectives of individuals interviewed.
Process tracing (Barnett and Munslow 2014; Davies, Laidlaw, and Rogers 2016; Talcott and Scholz 2015; White and Phillips 2012)	Uses qualitative methods to determine relative weight of evidence for causal links between activities and outcomes. The evidence is used to confirm whether mechanisms match predicted hypothesis. Comes from the analysis of historical events.	<ul style="list-style-type: none"> To see if results are consistent with the hypothesised mechanisms of action and to see if alternative explanations can be ruled out. Intervention needs to be at a relatively mature stage and some level of meaningful change has occurred. 	<ul style="list-style-type: none"> Make strong causal claims about what mechanism(s) caused a given set of outcomes in any given case. Requires sufficient time and human resources to enable participatory iterations of analysis and discussion with stakeholders.

(Continued)



Table 3. (Continued).

Method	Description	Use and timing	Strengths/Weaknesses/Considerations
Qualitative comparative analysis (QCA) (Befani 2013; Davies 2016a, Davies 2016b; Jordan et al. 2011; Kane et al. 2014; White and Phillips 2012)	A theory-driven approach used to examine the relationship of a priori outcomes of interests and the conditions hypothesised to influence the outcome. Qualitative data is converted to quantitative data (either binary or ordinal data) and tabulated for each condition and outcome. Patterns in the resulting data table are identified to highlight pathways of conditions that produce an outcome.	<ul style="list-style-type: none"> To test existing theories and new assumptions and formulate new theories. To understand the context under which interventions work and how different implementation strategies effect outcomes. Potential to support short cycle learning about the effectiveness of specific activities being implemented during a project's lifespan. However, quite a time consuming process 	<ul style="list-style-type: none"> Provides causal inference. Does not account for temporality. Can use relatively small and simple data sets. Strong external validity. Allows for the generalisation of findings from a relatively small number of cases and offers the ability to identify different pathways of condition combinations that lead to a similar outcome. Do not need to pre-specify causal pathways between activities and outcomes. May require more data as likely there will be a wider range of interventions and outcomes where relationships are possible. Does not require a baseline or comparison group. Does not provide an estimate of magnitude of effect. Quantitative coding of qualitative data speeds up data analysis. Findings presented in a dashboard, make them easy to interpret. Can identify unintended consequences Reflects perspectives of individuals interviewed; may not be generalisable to other settings. The QuIP incorporates features of a range of other qualitative approaches, including contribution analysis, process tracing, outcome harvesting and realist evaluation Aims to address the challenges of confirmation bias through blindfolding interviewers and respondents from knowing the full details of the intervention evaluated
Qualitative impact assessment protocol (QuIP) (Copestake 2014; Copestake, Morsink, and Remnant 2018b; Copestake et al. 2018a; Copestake and Remnant 2015)	Outcomes are explored with programme or intervention intended beneficiaries, to identify those factors beneficiaries perceive to be driving changes. Interviewers are blinded to the theory of change and project being assessed. Ask about causal drivers of change in selected areas of respondent's life. Data is coded quantitatively, highlighting whether reasons given for change confirm the hypothesised causal pathways. Code whether evidence is explicit (that is referenced project) or implicit.	<ul style="list-style-type: none"> Undertaken at a single point in time; although could be repeated to examine change over time Particularly useful where evaluation has not been incorporated from a programme's outset or where the context is highly changeable. Examines whether interfaceries on intended beneficiaries Provides both confirmatory (for example to test theory of change) and exploratory learning (for example open to unanticipated drivers and outcomes). 	<ul style="list-style-type: none"> Findings presented in a dashboard, make them easy to interpret. Can identify unintended consequences Reflects perspectives of individuals interviewed; may not be generalisable to other settings. The QuIP incorporates features of a range of other qualitative approaches, including contribution analysis, process tracing, outcome harvesting and realist evaluation Aims to address the challenges of confirmation bias through blindfolding interviewers and respondents from knowing the full details of the intervention evaluated

(Continued)

Table 3. (Continued).

Method	Description	Use and timing	Strengths/Weaknesses/Considerations
Rapid assessment process/Rapid assessment methodology (Beebe 2001; Butler 1995; Harris, Jerome, and Fawcett 1997; Hildebrand 1981; Manderson and Aaby 1992; Schünemann 2015; Vlassoff and Tanner 1992)	Highly focussed team based ethnographic approach; uses IDIs, FGDs and observations. Three major features: (1) a systems approach; (2) triangulation of data; (3) interactive data collection process to quickly develop a preliminary understanding of a situation from the insider's perspective.	<ul style="list-style-type: none"> • Could be undertaken at any stage of the programme. • Undertaken at a single point in time. • Aims to collect only relevant and necessary data; makes more rapid and cost-effective than traditional qualitative approaches. • Teams of interviewers may be used to rapidly collect information with the study completion expected within four to six weeks. 	<ul style="list-style-type: none"> • Ability to adjust investigations to reflect local conditions and specific situations. • Involve the community in both defining community needs and seeking possible solutions. • Adopts the principle of adequacy rather than scientific perfection. • Subject to both respondent (courtesy bias, social acceptability/political correctness bias, positional bias/attribution bias, self-serving bias and self-importance bias) and evaluator biases (contract renewal bias, friendship bias, and similar-person bias). • Provides more timely information than a systematic review by omitting stages of the systematic review process. • Less rigorous than a systematic review; search is not as comprehensive, may not double screen/extract, limited interpretation of findings, and so forth.
Rapid review/Expedited review, Rapid evidence summary (Ganann, Ciliska, and Thomas 2010; Grant and Booth 2009; HEARD Project 2018; Tricco, Langlois, and Straus 2017, Tricco et al. 2015)	A form of evidence synthesis. Methods vary; follows systematic review approach but places greater number of restrictions; for example fewer databases searched, time and setting restrictions or omits some processes to produce information in a timely manner.	<ul style="list-style-type: none"> • To identify new or emerging evidence on a topic, to assess what is already known about an intervention. 	<ul style="list-style-type: none"> • Assumes linear causal pathways. • Findings will be context specific.
Root cause analysis (Hubbard 2010; Peerally et al. 2017)	A method of structured risk identification and management. Not a single technique; a range of approaches and tools drawn from fields including human factors and safety science used to establish how and why an incident occurred in an attempt to identify how it, and similar problems, might be prevented from happening again. Analysis aims to establish a sequence of events to understand the relationships between contributory factors, the root cause and the defined problem.	<ul style="list-style-type: none"> • Typically undertaken to identify the cause after an adverse event has happened. • Can be used to forecast or predict 	
	Undertaken by a small team of stakeholders and facilitated by an expert.		

(Continued)

Table 3. (Continued).

Method	Description	Use and timing	Strengths/Weaknesses/Considerations
Statistical process control (Benneyan, Lloyd, and Plsek 2003; Fereday 2015)	Combines time series analysis methods with graphical presentation of data. Output or outcome data are plotted over time against statistical limits to identify if observed variation in an outcome deviates from the expected level of variations. Signals when the data deviates from predictions.	<ul style="list-style-type: none"> To determine whether changes in processes are making a difference to outcomes and/or to detect potential differences arising from different implementation strategies between sites. Undertaken continually throughout programme using data collected at standard intervals. 	<ul style="list-style-type: none"> Measures short-term outcomes. Limited measurement of longer-term impact. Requires ongoing data collection. Requires data collection, analysis and feedback to be completed as close to real time as possible Able to detect process changes and trends from an early stage in the programme; does not rely on reaching a pre-specified sample size – data limits adjusted when reason to believe current limits are not appropriate to provide adequate signals for action. Can change indicators or incorporate new outcomes overtime Potential bias from temporal trends Participants must be willing to engage in an honest and reflexive discussion. Findings will be context specific.
Strategy testing (Ladner 2015)	Participatory process for adapting theory of change over time. Initial theory of change represents best guess, which is examined on a regular basis to determine whether the assumptions are still valid.	<ul style="list-style-type: none"> To articulate and capture changes in the programme theory. A structured conversation undertaken with relevant stakeholders every 3–4 months throughout project. 	

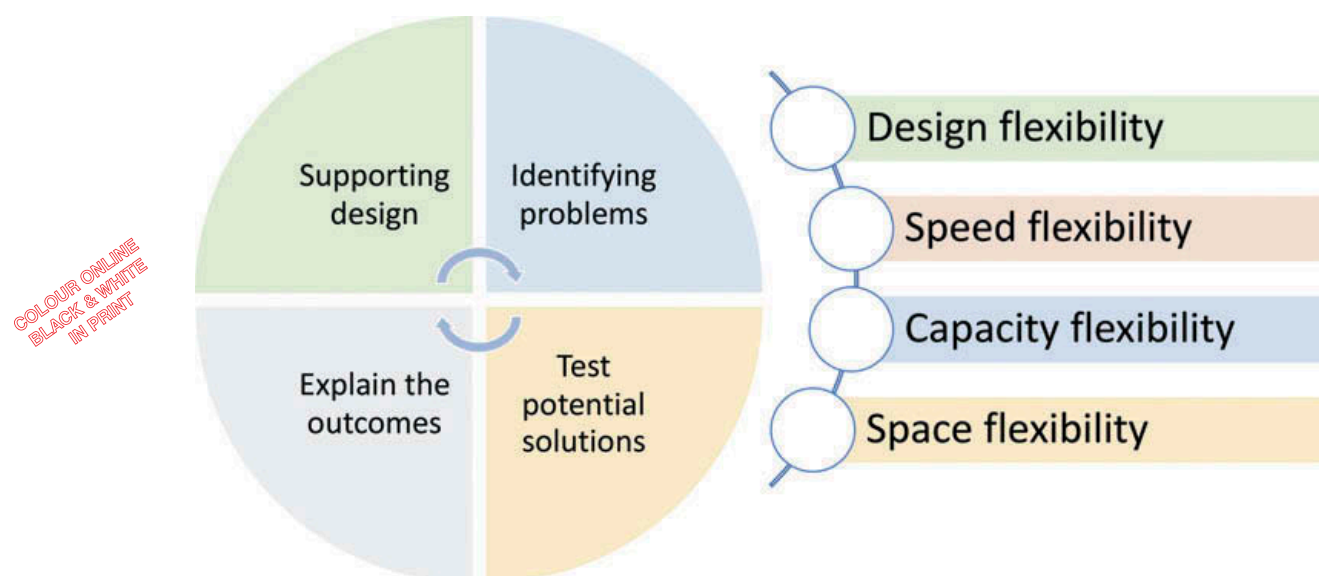


Figure 1. Framework for a timely approach to evaluation.

Design. The extent to which a method can respond to emerging insights and unexpected or unintended consequences once it has been designed, gained approvals, and its implementation is underway.

Speed. Ability of the method to adapt to time constraints and requirements. It considers the time required for design, data collection, analysis, reporting and feedback of data, and the potential to speed the process up. 160

Capacity. The level of skill required for design, data collection and analysis, and the extent to which there is flexibility around any of these.

Space. The ability of the method to adapt to different places and contexts. 165

Critique of methods against framework for application in international development

To illustrate the use of the analytical framework, we mapped a sub-set of methods against the four purposes and critiqued the applicability of the methods for a more timely approach using the four dimensions of timeliness and flexibility (Table 4). It is likely that over the course of a programme or intervention different methods will be needed to answer different evaluation questions and that the timescales and context will place restrictions on the suitability of different methods. A number of the methods identified can be used for multiple purposes and in general are not stand alone. We discuss the application of these methods for different purposes and discuss some of the challenges identified in critiquing the methods against the dimensions. The methods selected are intended to provide examples of the use of the framework to determine the applicability of a method, they are not intended to indicate exclusivity of these particular methods for timely evaluation in international development. 170 175

Support design

Both qualitative and quantitative methods can be used to support the development and/or refinement of an intervention or programme. Examples include rapid assessment process (RAP), a method of highly focussed ethnographic research, which draws on qualitative methods including in-depth interviews (IDIs), focus group discussions (FGDs) and observations (Beebe 2001), and A/B 180

Table 4. Appraisal of purpose categories of evaluation methods for application in timely evaluation.

Method	Purpose				Dimensions of timeliness and flexibility			
	Support design	Identify problems	Test solutions	Explain outcomes	Design	Speed	Capacity	Space
A/B testing	Primary		Secondary		Moderate; test a priori outcomes – outcomes can be changed for each cycle of testing based on emerging information.	Rapid; measures short-term outcomes. Depends on timeliness of routine/programme data.	Requires statistical expertise, but potentially programmes/packages could be developed to be operated with less expertise	Potentially adaptive to changing contexts.
Adaptive RCT	Secondary	Secondary	Primary		Limited; modifications are pre-planned before data is analysed based on pre-determined decision rules. Outcomes specified at outset are maintained throughout.	Moderate; Combine phases of a trial, reducing the time needed.	High level design and analysis expertise required	Not adaptive to changing contexts once started. Decision rules have to be pre-determined.
Bottle neck analysis		Primary			Moderate; causal pathways hypothesised prior to data collection. Does not capture unexpected/unanticipated outcomes. If analysis repeated causal pathways can be adapted to reflect changes to the programme. Limited; time built in between implementing steps to identify need for and make modifications	Moderate; depends on timeliness of routine/programme data. Slower if collecting primary data.	Moderate analytical expertise required	Can be adapted to include new/different hypothesised causal pathways
Modified stepped wedge trial	Secondary	Secondary	Primary			Moderate; implementation over defined phases of time	High level design and analysis expertise required	Not adaptive to changing contexts once started. Decision rules have to be pre-determined.
MSC	Secondary		Primary		High; can capture unexpected or unintended consequence	Slow; Takes time to collect stories, manage selection panels and feedback findings.	Interviewees can be trained relatively rapidly, and skills increases over time. Coding, analysis and interpretation requires skills & experience. Understanding of theory required.	Highly adaptive to different and changing contexts

(Continued)

Table 4. (Continued).

Method	Purpose			Dimensions of timeliness and flexibility		
	Support design	Identify problems	Test solutions	Explain outcomes	Design	
QuIP		Secondary		Primary	High; can capture unexpected or unintended consequence	<p>Speed</p> <p>Depends on time taken to collect data. Reduced analysis time by converting qualitative data into quantitative.</p> <p>Capacity</p> <p>Interviewees can be trained relatively rapidly, but skill increases over time. Coding, analysis and interpretation require skills and experience. Understanding of theory required.</p> <p>Space</p> <p>Adaptive to different and changing contexts</p>
RAP	Primary		Secondary	Secondary	High; grounded theory analysis allows shift in focus based on emerging findings. Inductive adaptation of interview guides	<p>Speed</p> <p>Rapid; estimated to be completed in 5–6 weeks.</p> <p>Capacity</p> <p>Interviewees can be trained relatively rapidly, but skill increases over time. Coding, analysis and interpretation require skills and experience. Understanding of theory required.</p> <p>Space</p> <p>Highly adaptive to different and changing contexts</p>
SPC		Primary	Secondary		Moderate; decision rules apply, however, outcomes can be expanded &/or adjusted over time.	<p>Speed</p> <p>Rapid; Does not rely on reaching pre-specified sample size so able to detect changes from an early stage. Depends on time to collect, analyse and feedback data. Can be very rapid where routine data is readily available.</p> <p>Capacity</p> <p>Visual output is easy to interpret. Computer programmes available to support analysis.</p> <p>Space</p> <p>Analyses highly adaptable across contexts, data recorded potentially difficult to change to include new indicators and data points recorded</p>

NOTE: Primary = main focus of the approach; secondary = possible, but not a main focus.

testing (also known as nimble RCTs, split tests, rapid-fire tests, bucket testing, **randomised** field experiments), a randomised trial in which participants are randomly assigned to receive a variation of the same intervention (Dibner-Dunlap and Rathore 2016; IPA 2016; Karlan 2017). 185

RAP is undertaken at a single point during the study to quickly develop a preliminary understanding of a situation. RAP was initially developed to support the evaluation of farming systems within a single planting season (Butler 1995; Hildebrand 1981) and has been used to develop interventions in health, for example to inform the development of tailored interventions for oral rehydration salts for diarrhoeal disease prevention within a limited time (Manderson and Aaby 1992) and for assessing operational challenges in the delivery of Long Lasting Insecticidal Nets (Theiss-Nyland et al. 2017). Qualitative methods, such as IDIs and FGDs are able to adapt to rapidly changing contexts or shifting priorities over time; inductive adaptation of interview guides and discussion themes on a daily basis can respond to emerging or unexpected findings. Transcription, translation, coding and analysis for in-depth exploration of the data are time consuming but in RAP for example, adaptations for rapid use are made that enable completion of a study within a relatively short time period. Teams of interviewers may be used to rapidly collect information with the study completion expected within four to six weeks (Harris, Jerome, and Fawcett 1997; Vlassoff and Tanner 1992). The emphasis is on adequacy of data for the purpose, rather than high level of precision. 190 195 200

RAP methods can be undertaken before programme implementation, when there is ambiguity about the scale and nature of the problem and what is needed to address a problem. It can be used to characterise the setting, assess whether a proposed programme or intervention addresses a particular need, is likely to be acceptable, and the feasibility of delivery, **and so forth**. The agility and speed with which RAP can be undertaken make it particularly useful when a problem has been identified to rapidly determine potential refinements to an intervention or programme and/or its delivery. Where differences in implementation have been identified then qualitative methods can explore reasons for 'positive deviance' to develop hypotheses about what has allowed the intervention or programme to succeed in some settings/participants when it has failed in the majority. Qualitative methods can be used to generate hypotheses about how a programme or intervention might work, particularly when, for example, a realist approach is taken and context-mechanism-outcome configurations developed (Manzano 2016; Pawson and Tilley 2004). This can usefully inform the design of future evaluation activities, including identifying relevant outcome measures. 205 210

Where there is a greater understanding of the type of intervention that is to be implemented methods, such as A/B testing can be used to refine the intervention before wider scale up and testing. A/B testing is most suited to testing small modifications to a programme's design or messaging, where the changes introduced are intended to result in immediate change (Optipedia n.d.). The focus on short-term outcomes, such as use and uptake, enables rapid testing of elements of a programme within a relatively short time frame but does not provide insight on longer-term impacts. As such A/B testing is particularly useful at the design or pilot stage of a programme and for answering questions about the early stages of a programme's theory of change. A/B testing has been used in South Africa to examine the impact of advertising content on demand for loans (Bertrand et al. 2010) and in Pakistan, Turkey, South Africa, Jordan, Bolivia, Peru and the Philippines to study the impact of varying message content of financial products in (Dibner-Dunlap and Rathore 2016; Karlan et al. 2016). To be most effective, A/B tests rely on good quality routine or administrative data and require a large sample size to be able to measure small incremental changes. 215 220 225

Identify problems

We illustrate two example of quantitative methods for identifying problems; statistical process control (SPC), which combines time series analysis with graphical presentation of data, and bottle neck analysis, which identifies blocks in the implementation process. Qualitative methods are also important in highlighting unintended or unanticipated consequence of existing interventions. 230

SPC originates from manufacturing and has been used for monitoring and quality improvement in healthcare. It is a statistical method that combines time series analysis methods with graphical presentation of data to identify if observed variation in an outcome deviates from the expected level of variations (Benneyan, Lloyd, and Plsek 2003; Fereday 2015). SPC is undertaken continually throughout a programme using data collected at standard intervals provided routine or operational data is available. It does not rely on reaching a pre-specified sample size as the statistical limits are varied accordingly; limits are adjusted when there is reason to believe that current limits are not appropriate to provide adequate signals for action. This means that SPC is able to detect process changes and trends from an early stage in the programme and that different outcome measures can be tracked overtime. The review did not identify examples of SPC having being used in a development context.

SPC is useful in situations where the context is complex and changeable as new outcomes can be dropped or added to the analysis as the intervention or programme is modified and its underpinning theory of change evolves. A highly adaptive approach to programming is likely to increase the number of outcome indicators that are measured. Changing outcomes is possible provided they are already available or easy to add to existing data collection tools. Where new data has to be collected this may have cost implications. SPC can also be used to detect potential differences arising from different implementation strategies between sites. This can highlight important differences that might warrant further investigation for example using qualitative methods to explore positive deviants.

Bottleneck analysis is one of three similar approaches to identifying the 'component(s) of a system that limits the overall performance or capacity' (O'Connell and Sharkey 2013; Rio et al. 2015). Two related ideas are cascade analysis and community or systems effectiveness (Dellicour et al. 2016; Garnett et al. 2016; Webster et al. 2013). In each case a number of steps that link the population intended to benefit from an intervention and the population that do benefit are identified and assessed. Each step is conditional on the previous one having been met and only the population left at the end of all the steps would be anticipated to have achieved the desired outcome. The relative size of the population lost at each step might indicate where the most urgent action is needed. For example, a bottleneck analysis of maternal and newborn health interventions in rural areas of the United Republic of Tanzania, found the largest bottleneck in one region was the availability of equipment, drugs and human resources in the facility, while in another the largest bottleneck was clinical practice (Baker et al. 2015). These methods are usefully combined with qualitative approaches to explore why the bottleneck has occurred and identify potential modifications to a programme.

Bottleneck analysis assumes a linear process; that achieving one step is a necessary condition to achieving the next. This implies that the hypothesised theory of change is the only route through which change can occur. To assess if this assumptions holds, requires an understanding of whether the population in one stage is the same as the population in the next, to ascertain whether it is a 'necessary' condition or whether other steps, not captured in the theory of change, might be sufficient to achieve the desired change (Davies 2014). The analysis could be adapted to reflect changes in understanding of necessary and sufficient conditions and as the programme's theory of change evolves, provided data is available on the relevant outcomes.

Such analyses are often undertaken at a single point in time and provide a snap shot of need. Where routine or programme data is available the analysis can be undertaken relatively rapidly and could be repeated to assess whether the bottlenecks identified and size change overtime.

Test potential solutions

Experimental methods are used to assess the effectiveness of interventions or programmes and to ascertain causal relationships. Recent innovations including adaptive randomised control trials (RCTs) and modified stepped wedge trials present real opportunities for these methods to usefully support timely approach to evaluation. Their use for complex interventions in international

development however, has been highly restricted to date. The review identified one protocol for an adaptive RCT and one protocol for a modified stepped wedge trial in international development settings (Choko et al. 2017; Wechsberg et al. 2017).

Adaptive RCTs can be used to test multiple interventions in parallel before applying stopping rules as the evidence stacks up. This method may be particularly useful where it is not clear which interventions are most likely to be effective to achieve similar outcomes. The design includes multiple rounds of interim analysis that allows interventions that are not performing according to predetermined criteria to be terminated (Bothwell et al. 2018; Kairalla et al. 2012; Mahajan and Gupta 2010). In addition to starting or stopping interventions modifications can include: adjusting the study population and sample size; and outcome-adaptive randomisation in which treatment allocation is skewed to those treatments that appear to be doing better. Potential modifications, and the criteria for implementing changes, need to be pre-specified based on decision rules in the study protocol.

The inclusion of a period of 'reflection' between each step of implementation in a modified stepped wedge trial makes this method useful where the basic form of an intervention has been decided upon at the outset but enables testing of the acceptability, feasibility and effectiveness of the intervention as it is implemented. Between steps formative research, including surveys, IDIs and FGDs, assess the acceptability and feasibility of implementing the intervention or programme and, where relevant, identify a revised plan to be implemented in the next step. At the end of the study, it would be possible to compare the effect of the overall package of interventions on the pre-specified outcomes as in the original study, but additionally provides an evidence-based refined delivery plan for roll-out in other areas.

Both methods can be combined with methods, such as SPC to determine whether causal mechanisms are being activated as anticipated, as well as qualitative methods to understand the mechanism by which an intervention has impact, capture unanticipated outcomes and/or the influence of context (Stetler et al. 2006). The value of adapted or modified trials lies in their ability to make adjustments to the intervention or trial design as data is being collected, without undermining the validity or integrity of the study (Bhatt and Mehta 2016; Bothwell et al. 2018; Kairalla et al. 2012; Korn and Freidlin 2017; Lang 2011; Thorlund et al. 2018; Villar, Bowden, and Wason 2017). This provides both ongoing learning during the programme and confirmatory learning at the end of the trial, which could be generalised to other settings. Such designs require significant investment and expertise can increase trial complexity and require sophisticated statistical techniques for the analysis.

Explain outcomes

Explaining outcomes draws primarily on qualitative methods to gather stakeholder and beneficiaries' perceptions of interventions and programmes or elucidation of their causal mechanisms. Examples include most significant change (MSC) and qualitative impact assessment protocol (QuIP). Both methods are undertaken retrospectively when sufficient time is anticipated to have passed to warrant examination of impact of an intervention or programme. The methods start by assessing whether meaningful change has occurred and work backwards to determine whether change can be attributed to the specific intervention (Beach and Pedersen 2013; Lacouture et al. 2015).

MSC was originally developed as a form of participatory impact monitoring (Davies 1996), to be used in a decentralised and participatory rural development programme, where standardised pre-defined indicators would not work. In each reporting period (initially 3 months), programme participants were asked to identify what they thought was the MSC, and its consequences. Stakeholder panels review these stories to identify the most significant and the consequences for the NGO's future work. In the decades since then MSC has been used in a wide variety of programmes for both evaluation and monitoring purposes. Many different selection structures have been designed to fit the different kinds of programmes and stakeholders involved (Davies and Dart 2005). MSC is particularly valuable in highly complex settings where it is not known which activities are likely to have led to change and where causal mechanisms have either not been articulated at the project outset or cannot be agreed upon between stakeholders.

QulP assesses impact through narrative causal statements from programme or intervention intended beneficiaries. The QulP takes on the challenge of achieving sufficient credibility using timely qualitative methods in a way that can be both confirmatory (testing a theory of change) and exploratory (open to the unanticipated drivers and outcomes) (Copestake 2014). It was developed through a grant to evaluate rural livelihood adaptation projects in Malawi and Ethiopia but has since been used to conduct relatively rapid studies in many other fields, including assessment of the social impact of ongoing programmes to promote decent work in Mexico, community self-organisation in Uganda and improved housing in India. (Copestake, Morsink, and Remnant 2018b; Copestake and Remnant 2015). The QulP incorporates features of a range of other qualitative approaches, including contribution analysis, process tracing, outcome harvesting and realist evaluation. It builds on ongoing quantitative monitoring of key indicators using semi-structured interviews and FGDs. Its potential as a timely and flexible approach is enhanced by requiring neither a baseline nor a comparison group. But like other forms of contribution analysis it tests the existence of causal pathways, but does not generate estimates of the magnitude of causal effects. Field data generated on drivers of change is open-ended and exploratory, because the field team is deliberately not informed of project theory (or even the identity of the project being evaluated). But a critical part of the job of the analyst is to code the drivers of change identified according to whether they do explicitly or implicitly align with project theory or not. The QulP aims to address the challenges of confirmation bias (where what people say is framed by how they are interviewed and possibly influenced by what they think you want to hear) through 'blindfolding' interviewers and respondents from knowing the full details of the intervention evaluated (Copestake et al. 2018a).

These methods are generally undertaken at a single point in time, although they can be repeated to examine how perspectives change over time; in this way these method can assess both short and longer-term outcomes and can provide insights into whether a programme is having its intended impact and which activities are responsible for any observed change. These methods are particularly valuable where the interventions being implemented or the context are highly complex and changeable. They are also valuable where evaluation has not been incorporated from a programme's outset.

Both methods have the potential to be used for hypothesis testing, they examine what was achieved and how, to understand the relative importance of different activities undertaken. However, there is considerable flexibility as data collection is not restricted to pre-specified outcomes. This allows evaluators to capture unexpected outcomes and mechanisms of action, and can lead to new hypotheses and theories being generated. The timeliness of evidence can also be enhanced (relative to more traditional methods of qualitative research) by adopting more structured protocols for data coding, analysis and visualisation. The QulP method has sought to speed up the process of synthesis and reporting by speeding up data analysis and reporting through use of bespoke spreadsheets, and interactive dashboards to supplement more formal reports.

Discussion

We set out to develop a framework to identify, categorise and critically appraise methods that can support a more timely approach to evaluation of international development programmes. We identified both quantitative and qualitative methods that can be used for different purposes, namely: supporting design, identifying problems and testing and explaining solutions. We suggest methods are selected based upon the purpose of the evaluation. This analysis highlights that different methods can fulfil multiple purpose; the particular method to be used should be selected based on the specific time-needs and flexibility of the programme.

Our review found there to be a dearth of examples of the application of methods being explicitly used for more timely approaches to evaluating international development programmes. Reasons for this may include that those conducting such evaluations rarely disseminate their findings through

peer reviewed publications or through widely accessible grey literature. We are optimistic that there is significant potential for timely evaluation to improve international development outcomes. Realising this challenge however will require further understanding of a number of core issues and further work to develop and test methods to be used for timely evaluations. We reflect on some key issues that were repeatedly raised in discussions and in the literature.

385

To detect change in a timely manner relies on the analysis of outputs and short-term outcomes to indicate change rather than longer-term impacts. This particularly applies to quantitative methods, such as SPC, A/B testing and interim-analysis of adaptive or modified trials. The use of shorter-term outcomes run the risk of falsely detecting treatment effects or prematurely discarding promising interventions that do not show an impact at an early stage. It is therefore important to recognise the short time horizon of applicability of the findings and conclusions drawn need to be viewed with caution as assessing impact over a longer period might lead to different conclusions or other information emerging as causal processes work over different time scales (Woolcock 2009).

390

The advantage of methods like adaptive and modified trials is that they can also provide confirmatory learning at the end of the trial, demonstrating whether an intervention had the intended impact by measuring pre-defined outcomes over the entire course of the trial. Outcomes are selected based on hypothesised causal chains. These methods should be combined with qualitative methods to pick up unanticipated outcomes. When using methods, such as SPC that have the flexibility to change the outcomes measured overtime, researchers should consider the value of including some constant or 'bedrock' indicators that don't change over the life of the programme to support an understanding of the longer term impact of projects (Barr 2015).

395

We did not identify any documentation of the impact that measuring and basing decision on shorter term outcomes has in this setting through the literature review. However, during the symposium concerns were raised that these approaches might cause researchers to become too focused on short term outcomes at the expense of the longer term impacts and the impact on rigour. More research is needed to understand the validity and rigour of using more timely methods compared to endline analysis. This could be tested for example in a trial with different forms of timely evaluation as the different arms, for example different timings of feeding back results, with different data sources informing the results.

405

Using pre-existing data can reduce the time and resources needed for quantitative methods. However, many development programmes have weak monitoring systems which make them less likely to be easily evaluable. Timeliness for many of the methods will therefore depend on the ability to collect, process and analyse data in a timely fashion. The challenge is to better leverage time series data from service delivery platforms and to make such data useful (that is captures relevant outcome indicators in a timely manner) and of sufficient quality (that is measures needed to enhance completeness and accuracy of data).

410

415

The ability of routine data to respond to shifting priorities over time and the amount of time required for data collection and analysis is variable depending upon the scale and ownership of the data collection system. While changes to the indicators in national-level routine systems are a major undertaking, other forms of routine data capture, such as programme monitoring data, may be more flexible and outcomes measured could be adjusted over time. The key therefore is in the initial design and whether an expectation of the need for flexibility has been built into the system. Where high quality routine data is available, then analysis is generally very rapid.

420

In settings where routine data is not available, innovative approaches to accessing routine data offer real potential (DFID 2012). For example, the American Refugee Committee uses digital technology to collect highly focussed satisfaction data from refugees in camps in Uganda, Rwanda, Somalia and Sudan (Peters 2018). While, during the 2013–16 Ebola outbreak in West Africa real-time data surveys were undertaken resulting in significant lessons learned on the rapid collection, coordination and use of large amounts of data using new technologies and on coordination of this data amongst partners (Cori et al. 2017). The analysis of big data is already common place in the private sector; used for consumer profiling, personalised services and predictive

425

430

analysis being used for advertising (UN Global Pulse 2012). Technology that offers increasing opportunities for real time data analytics and their application should be explored more in development programmes.

The general consensus from the public symposium and literature review was that the use of mixed methods should be encouraged; quantitative approaches should be complemented for their interpretation, by process data, which is often qualitative. Mixing of methods can ensure a greater sensitivity amongst evaluators towards the potential threats to the validity of conclusions (Ton 2012). It has become a general expectation that impact evaluations be accompanied by a process evaluation and a similar approach makes perfect sense when considering timely evaluation within an ongoing programme.

A mixed methods approach may involve using complementary methods of data collection, but may also mean mixing or combining of theories, hypotheses, analyses and conceptual or analytical frameworks (Bamberger 2012). Innovative approaches to mixing methods, stemming from the field of political science, have recently been proposed. Goertz's 'research triad' is a multi-methods approach which links not just quantitative (cross-case) with qualitative (within-case) inferences, but adds a third approach of the elucidation of causal mechanisms through for example, process tracing (Goertz and Mahoney 2012). Amongst the qualitative approaches the interpretative approaches tend to have a focus on for example, the influence of power and the meaning behaviours, whilst a subset of methods is concerned with causal inference, mechanisms and generalisation (Goertz and Mahoney 2012).

Stakeholder engagement is essential to ensure efficient incorporation of learning from timely evaluation into programme adaptations that can successfully be implemented. This can increase the utility of an evaluation to support programme improvement – an approach espoused by Patton called 'utilisation focused evaluation' (Patton 2008), in which end-users are identified and engaged from outset to guide other decisions that are made about the evaluation process. This has great benefits, though it also requires sufficient time and resources, as well as willingness on the part of the stakeholders. Evaluation also needs to be responsive such that results are available whilst there is momentum and engagement amongst staff. Sometimes staff may have solved problems that the evaluation later highlights the presence of, and therefore the evaluation is no longer relevant for pushing programme improvement.

The programmes within which the timely evaluation framework and approaches are applied

There is a close link between what the evaluation methods are trying to do, and the ability of programmes to incorporate and act on what they tell us either at programme outset, through adaptations over time that are responsive to monitoring data, or in acting on the results of comparative or explanatory studies on programme options or performance. A central issue to these are the intersection between programming flexibility/adaptability and the timing with which data from evaluation is 'received' and how this links to programming cycles.

It was argued at the symposium that programme improvement is only really possible when: (1) programmes are small; (2) there is a specific intention to learn and adapt; (3) when results are immediately available; (4) when changes to the programme are small-scale within the capacity of the programme to deliver; and (5) when programmes have time to try out various options before rolling out to reach a large number of beneficiaries (Aly Visram personal communication). Large scale improvements are difficult if not impossible to implement, especially because they require significant investment. Large scale improvements are also likely to be beyond the financial capacity of programmes that have pre-budgeted based on a fixed plan of action. The proposition of achievement through small incremental changes is supported by the idea from evolutionary theory of 'the adjacent possible' (Srivastava 2014).

Effective use of data requires appropriate data, that reaches the right people, who understand the data as presented, are able to transform it as required, and have the power to make decisions or have access to those who do. The guidance on change must then be produced and transferred back to implementers who are able, and willing, to put changes into action. The presence of programme and institutional structures required to support this process, which in itself is complex, will vary. 480

Uncertainty over what evidence might be needed and when is often compounded by delays in the time it takes commissioners and evaluators to respond. Empirical evidence on the processes involved in generating evidence is lacking, partly perhaps because the scope for generalising usefully about it is limited by context-specificity. Having set out to develop a more agile approach to collecting 'good enough' evidence in the form of the 'QuIP' James Copestake reflected at the public Symposium, on practical obstacles to doing so. 485 490

Starting with the demand side, delays arise in securing agreement on the design, budget, release of sample-frame data, clarity on the theory of change needed to guide data coding and on obtaining ethics approval sometimes across more than one institution. These are particularly likely when the commissioner seeking an evaluation and the organisation executing the activity being evaluated are distrustful of each other. Delays arise from variation in the nature of the primary intended audience and their expectation of what evidence should look like, which may range from a flexible data dashboard to a glossy report. The more controversial the findings (and hence perhaps the more important), the more the likelihood of lengthy negotiation over an 'acceptable' final draft. Meanwhile, on the supply side, the challenge of mobilising appropriate and available staff for data collection is often compounded by problems securing permission to enter the field, finalising contracts and securing ethical approval (Gamble 2006; Patton 2013; Portela et al. 2015). 495 500

There is a need to test the scope of timely evaluation methods and to determine which programmes they can or should be applied to. There is limited evidence in particular for outcome evaluation methods presented here (adapted RCTs and modified stepped wedge trials), which might support large scale testing and change. 505

Assessing the impact of timely evaluation

Timely evaluation approaches are likely to be more time and resource intensive. All of the methods presented are likely to be resource intensive and require more data to be collected than traditional evaluation methods. Methods that do not test a specific causal mechanism need to capture a wider range of outcomes and causal pathways. Whilst, methods that aim to rapidly test changes or compare multiple-interventions rely on ongoing or repeat measurement of data. The methods are anticipated to represent overall value for money as they result in the programme having a higher chance of success. However, the impact/benefits of undertaking more timely approaches to evaluation are not well understood (O'Donnell 2016). There is therefore a need to determine whether undertaking a timely evaluation does lead to greater impact than traditional approaches and represent value for money. 510 515

It is important to understand the implications of learning more for this time on our ability to learn more for next time. Where an intervention changes over time there is a need to identify when it becomes an entirely new intervention and to recognise when the use of these methods become an intervention in themselves (Portela et al. 2015). If this is the case the use of these methods may need to be incorporated into interventions being replicated in different settings. It is questionable then whether we can learn anything on scaling up or replication in other settings using these approaches. It is necessary to understand the nature of implementation and the degree to which evaluation activities influence and contribute to the overall results of a programme. 520

Limitations of our approach

525

There were several limitations to our approach. Our scoping seminar and public symposium were interesting and exciting events, which provided an opportunity for broad discussion of timely evaluation within international development. Although in setting the agenda and selecting speakers we attempted to focus some of the discussions, the topic was new for many participants and therefore the discussions quite broad.

530

Reviewing the literature on this topic proved to be extraordinary difficult due to the wide range of terminologies around timely evaluation, programme improvement and adaptive learning. Many of the methods we identified were specific to certain niches for example, quality improvement initiatives. There were also a range of terminologies for what in effect were very similar methods. In addition to problems in terminology, there were many examples of methods being advocated for and described without any examples of their practical application or critique of this application.

535

Although we attempted to embrace a wide range of sectors in our paper, the experience of the majority of the author team, and participants of the scoping session and public symposium is in the health sector and therefore most of our examples are from the health sector. We hope however, that our framework and discussion of approaches and methods will provide a starting point, which can be applied across sectors.

540

Identification, categorisation and better selection of methods for timely evaluation within specific programmes can only go so far in improving outcomes: uncertainty will always remain about 'what works, for whom and under what circumstances'. Borrowing this mantra from the tradition of realist evaluation is not an accident because a complexity ontology is what underpins it, and its recognition that evaluation is unavoidably political, as well as technical (Pawson 2013).

545

Recommendations for further research

Based on our discussions and review of the literature we recommend further research on timely evaluation including:

Testing and development of framework

550

The framework should be tested to ensure fit for purpose. Workshops convening relevant stakeholders including researchers, implementers and decision makers could assess the utility of the framework for selecting methods and determining the optimum mix of methods for addressing different development projects being conducted in different contexts and settings. Through testing would also identify research priorities for developing new or adapting existing methods to meet the needs of a more timely approach to evaluation.

555

Developing guidelines and best practices

The framework should be developed further to provide guidance on best practices on timely evaluation for programme improvement for different types of projects within different contexts. This would involve formulating a matrix of recommended methods with guidance on their applicability for different projects, contexts and sectors, for example education and agriculture.

560

Evaluating adaptive management interventions

While the flexible approaches underlying adaptive management are very promising, these remain to be rigorously evaluated.

Conducting adaptive trials

565

The application of adaptive trials to multi-component interventions where different packages of configurations is tested, where there are ethical issues and decisions have to be made quickly. For example, humanitarian assistance interventions would be one of such cases.

Conclusion

There is significant potential for more timely evaluation to improve international development outcomes. Despite the availability of new approaches to mixing and adapting existing methods and the potential for new technologies to enhance data collection, there is a dearth of examples of their application. Research is needed to provide an empirical evidence base upon which to further develop and appraise the application of these methods, across sectors and contexts within international development.

Disclosure statement

AQ4 No potential conflict of interest was reported by the authors.

Funding

AQ3 This work was supported by the Department for International Development [203569].

References

- Andrews, M. 2015. "Explaining Positive Deviance in Public Sector Reforms in Development." *World development* 74: 197–208. doi:10.1016/j.worlddev.2015.04.017.
- Baker, U., S. Peterson, T. Marchant, G. Mbaruku, S. Temu, F. Manzi, and C. Hanson. 2015. "Identifying Implementation Bottlenecks for Maternal and Newborn Health Interventions in Rural Districts of the United Republic of Tanzania." *Bulletin of the World Health Organization* 93: 380–389. doi:10.2471/BLT.14.141879.
- Bamberger, M. 2012. "Introduction to Mixed Methods in Impact Evaluation." *Impact Evaluation Notes* 3: 1–38.
- Barder, O., and B. Ramalingam. 2012. *Complexity, Adaptation, and Results*. Center For Global Development.
- Barnett, C., and T. Munslow. 2014. "Process Tracing: The Potential and Pitfalls for Impact Evaluation in International Development." Summary of a Workshop held on 7 May 2014 (No. IDS Evidence Report 102). Institute of Development Studies.
- Barr, J. 2015. *Monitoring and Evaluating Flexible and Adaptive Programming*. Itad.
- Beach, D., and R. B. Pedersen. 2013. *Process-Tracing Methods: Foundations and Guidelines*. University of Michigan Press.
- Beebe, J. 2001. *Rapid Assessment Process: An Introduction*. Rowman Altamira.
- Befani, B. 2013. "Between Complexity and Generalization: Addressing Evaluation Challenges with QCA." *Evaluation* 19: 269–283. doi:10.1177/1474022213493839.
- Befani, B., and J. Mayne. 2014. "Process Tracing and Contribution Analysis: A Combined Approach to Generative Causal Inference for Impact Evaluation." *IDS bulletin* 45: 17–36. doi:10.1111/1759-5436.12110.
- Benneyan, J. C., R. C. Lloyd, and P. E. Plsek. 2003. "Statistical Process Control as a Tool for Research and Healthcare Improvement." *BMJ quality & safety* 12: 458–464. doi:10.1136/qhc.12.6.458.
- Bertrand, M., D. Karlan, S. Mullainathan, E. Shafir, and J. Zinman. 2010. "What's Advertising Content Worth? Evidence from a Consumer Credit Marketing Field Experiment." *The quarterly journal of economics* 125: 263–306. doi:10.1162/qjec.2010.125.1.263.
- Bhatt, D. L., and C. Mehta. 2016. "Adaptive Designs for Clinical Trials." *The New England journal of medicine* 375: 65–74. doi:10.1056/NEJMr1510061.
- Biglan, A., D. Ary, and A. C. Wagenaar. 2000. "The Value of Interrupted Time-Series Experiments for Community Intervention Research." *Prevention science : the official journal of the Society for Prevention Research* 1: 31–49. doi:10.1023/A:1010024016308.
- Bothwell, L. E., J. Avorn, N. F. Khan, and A. S. Kesselheim. 2018. "Adaptive Design Clinical Trials: A Review of the Literature and ClinicalTrials.gov." *BMJ open* 8: e018320. doi:10.1136/bmjopen-2017-018320.
- Burke, L. E., S. Shiffman, E. Music, M. A. Styn, A. Kriska, A. Smailagic, D. Siewiorek, et al. 2017. "Ecological Momentary Assessment in Behavioral Research: Addressing Technological and Human Participant Challenges." *Journal of medical Internet research* 19: e77. doi:10.2196/jmir.7138.
- Busza, J., S. Teferra, S. Omer, and C. Zimmerman. 2017. "Learning from Returnee Ethiopian Migrant Domestic Workers: A Qualitative Assessment to Reduce the Risk of Human Trafficking." *Globalization and health* 13. doi:10.1186/s12992-017-0293-x.
- Butler, L. M. 1995. *The Sondeo: A Rapid Reconnaissance Approach for Situational Assessment*.
- Cellamare, M., S. Ventz, E. Baudin, C. D. Mitnick, and L. Trippa. 2017. "A Bayesian Response-Adaptive Trial in Tuberculosis: The endTB Trial." *Clinical trials (London, England)* 14: 17–28. doi:10.1177/1740774516665090.

- Choko, A. T., K. Fielding, N. Stallard, H. Maheswaran, A. Lepine, N. Desmond, M. K. Kumwenda, and E. L. Corbett. 2017. "Investigating Interventions to Increase Uptake of HIV Testing and Linkage into Care or Prevention for Male Partners of Pregnant Women in Antenatal Clinics in Blantyre, Malawi: Study Protocol for a Cluster Randomised Trial." *Trials* 18. doi:10.1186/s13063-017-2093-2. 620
- Connors, S. C., S. Nyaude, A. Challender, E. Aagaard, C. Velez, and J. Hakim. 2017. "Evaluating the Impact of the Medical Education Partnership Initiative at the University of Zimbabwe College of Health Sciences Using the Most Significant Change Technique." *Academic medicine : Journal of the Association of American Medical Colleges* 92: 1264–1268. doi:10.1097/ACM.0000000000001519. 625
- Copestake, J. 2014. "Credible Impact Evaluation in Complex Contexts: Confirmatory and Exploratory Approaches." *Evaluation* 20: 412–427. doi:10.1177/1356389014550559.
- Copestake, J., C. Allan, B. W. Van, Belay, M. Goshu, T. Mvula, P. Remnant, F. Thomas, and E. Zerahun. 2018a. "Managing Relationships in Qualitative Impact Evaluation of International Development: QulP Choreography as a Case Study." *Evaluation* 24: 169–184. doi:10.1177/1356389018763243. 630
- Copestake, J., M. Morsink, and F. Remnant, Eds. 2018b. *Attributing Development Impact: The QulP Case Book*. Rugby: Practical Action.
- Copestake, J., and F. Remnant. 2015. *Assessing Rural Transformations: Piloting a Qualitative Impact Protocol in Malawi and Ethiopia*. In: *Mixed Methods Research in Poverty and Vulnerability*. 119–148. London: Palgrave Macmillan. doi: 10.1057/9781137452511_6. 635
- Cori, A., C. A. Donnelly, I. Dorigatti, N. M. Ferguson, C. Fraser, T. Garske, T. Jombart, et al. 2017. "Key Data for Outbreak Evaluation: Building on the Ebola Experience. Philos." *Philosophical transactions of the Royal Society of London. Series B, Biological sciences* 372. doi:10.1098/rstb.2016.0371.
- Davies, R. 1996. *An Evolutionary Approach to Facilitating Organisational Learning: An Experiment by the Christian Commission for Development in Bangladesh*. University of Swansea, Centre for Development Studies. 640
- Davies, R. 2014. *Thinking about Set Relationships within Monitoring Data*. Rick on the Road.
- Davies, R., 2016a. "Qualitative Comparative Analysis [WWW Document]." *Better Evaluation*. Accessed 1 March 2018. http://www.betterevaluation.org/en/evaluation-options/qualitative_comparative_analysis
- Davies, R. 2016b. *Evaluating the Impact of Flexible Development Interventions Using a 'Loose' Theory of Change Reflections on the Australia-Mekong NGO Engagement Platform (A Method Lab Publication)*. London: Overseas Development Institute. 645
- Davies, R., and J. Dart. 2005. *The "Most Significant Change" (MSC). Technique: A Guide to Its Use*.
- Davies, R., J. Laidlaw, and P. Rogers. 2016. *Process Tracing [WWW Document]*. Better Evaluation.
- DDD, 2014. "Doing Development Differently [WWW Document]." *Doing Development Differently*. URL Accessed 17 February 2018 <http://doingdevelopmentdifferently.com/the-ddd-manifesto/> 650
- Dellicour, S., J. Hill, J. Bruce, P. Ouma, D. Marwanga, P. Otieno, M. Desai, M. J. Hamel, S. Kariuki, and J. Webster. 2016. "Effectiveness of the Delivery of Interventions to Prevent Malaria in Pregnancy in Kenya." *Malaria journal* 15: 221. doi:10.1186/s12936-016-1261-2.
- DFID. 2012. *Results in Fragile and Conflict-Affected States and Situations: How to Note*. Department for International Development. 655
- Dibner-Dunlap, A., and Y. Rathore. 2016. "Beyond RCTs: How Rapid-Fire Testing Can Build Better Financial Products [WWW Document]." *Innovations for Poverty Action*. URL Accessed 2 January 2018. <https://www.poverty-action.org/blog/beyond-rcts-how-rapid-fire-testing-can-build-better-financial-products>
- Earl, S., F. Carden, and T. Smutylo. 2001. *Outcome Mapping: Building Learning and Reflection into Development Programs*. Ottawa, Canada: International Development Research Centre. 660
- Eirich, F., and A. Morrison. n.d. *Guide 6: Contribution Analysis, Social Science Methods Series*. Scottish Government.
- Fereday, S. 2015. *A Guide to Quality Improvement Methods*. Healthcare Quality Improvement Partnership.
- Gamble, J. 2006. *A Developmental Evaluation Primer*. Canada: J.W. McConnell Family Foundation.
- Ganann, R., D. Ciliska, and H. Thomas. 2010. "Expediting Systematic Reviews: Methods and Implications of Rapid Reviews." *Implementation Science* 5: 56. doi:10.1186/1748-5908-5-56. 665
- Garnett, G. P., T. B. Hallett, A. Takaruzza, J. Hargreaves, R. Rhead, M. Warren, C. Nyamukapa, and S. Gregson. 2016. "Providing a Conceptual Framework for HIV Prevention Cascades and Assessing Feasibility of Empirical Measurement with Data from East Zimbabwe: A Case Study." *Lancet HIV* 3: e297–e306. doi:10.1016/S2352-3018(16)30039-X.
- Goertz, G., and J. Mahoney. 2012. *A Tale of Two Cultures: Qualitative and Quantitative Research in the Social Sciences*. Princeton University Press. 670
- Grant, M. J., and A. Booth. 2009. "A Typology of Reviews: An Analysis of 14 Review Types and Associated Methodologies." *Health Information & Libraries Journal* 26: 91–108. doi:10.1111/j.1471-1842.2009.00848.x.
- Green, D. 2015. *Doing Development Differently: A Great Discussion on Adaptive Management (No, Really)*. From Poverty to Power. 675
- Harris, K. J., N. W. Jerome, and S. B. Fawcett. 1997. "Rapid Assessment Procedures: A Review and Critique." *Human organization* 56: 375–378. doi:10.17730/humo.56.3.w525025611458003.
- HEARD Project, 2018. "Rapid Review Vs. Systematic Review: What are the Differences? [WWW Document]." *USAID*. Accessed 26 October 2018. <https://www.heardproject.org/news/rapid-review-vs-systematic-review-what-are-the-differences/>

AQ11
AQ12AQ13
AQ14

AQ15

AQ16
AQ17

AQ18

AQ19

- AQ20** UN Global Pulse 2012. *Big Data for Development: Challenges and Opportunities*. United Nations. 680
- Hildebrand, P. E. 1981. "Combining Disciplines in Rapid Appraisal: The Sondeo Approach." *Agricultural Administration* 8: 423–432. doi:10.1016/0309-586X(81)90037-6.
- Ho, L. S., G. Labrecque, I. Batonon, V. Salsi, and R. Ratnayake. 2015. "Effects of a Community Scorecard on Improving the Local Health System in Eastern Democratic Republic of Congo: Qualitative Evidence Using the Most Significant Change Technique." *Conflict and health* 9: 27. doi:10.1186/s13031-015-0055-4. 685
- AQ21** Hubbard, B. 2010. *Root Cause Analysis (Overview)*. Lean Learning Revolution!
- AQ22** IPA. 2016. *Introduction to Rapid-Fire Operational Testing for Social Programs (Goldilocks Deep Dive)*. Innovations for Poverty Action.
- AQ23** Jones, H., and S. Hearn. 2009. *Outcome Mapping: A Realistic Alternative for Planning, Monitoring and Evaluation (Working and Discussion Paper)*. Overseas Development Institute. 690
- Jordan, E., M. E. Gross, A. N. Javernick-Will, and M. J. Garvin. 2011. "Use and Misuse of Qualitative Comparative Analysis." *Construction Management and Economics* 29: 1159–1173. doi:10.1080/01446193.2011.640339.
- Kairalla, J. A., C. S. Coffey, M. A. Thomann, and K. E. Muller. 2012. "Adaptive Trial Designs: A Review of Barriers and Opportunities." *Trials* 13: 145. doi:10.1186/1745-6215-13-145.
- Kane, H., M. A. Lewis, P. A. Williams, and L. C. Kahwati. 2014. "Using Qualitative Comparative Analysis to Understand and Quantify Translation and Implementation." *Translational behavioral medicine* 4: 201–208. doi:10.1007/s13142-014-0251-6. 695
- AQ24** Karlan, D. 2017. *Nimble RCTs. A Powerful Methodology in the Program Design Toolbox*. Innovations for Poverty Action, Yale University.
- Karlan, D., M. McConnell, S. Mullainathan, and J. Zinman. 2016. "Getting to the Top of Mind: How Reminders Increase Saving." *Management science* 62: 3393–3411. doi:10.1287/mnsc.2015.2296. 700
- Kontopantelis, E., T. Doran, D. A. Springate, I. Buchan, and D. Reeves. 2015. "Regression Based Quasi-Experimental Approach When Randomisation Is Not an Option: Interrupted Time Series Analysis." *BMJ (Clinical research ed.)* 350: h2750. doi:10.1136/bmj.h2750.
- Korn, E. L., and B. Freidlin. 2017. "Adaptive Clinical Trials: Advantages and Disadvantages of Various Adaptive Design Elements." *Journal of the National Cancer Institute* 109. doi:10.1093/jnci/djx013. 705
- Lacouture, A., E. Breton, A. Guichard, and V. Ridde. 2015. "The Concept of Mechanism from a Realist Approach: A Scoping Review to Facilitate Its Operationalization in Public Health Program Evaluation." *Implementation science* : IS 10: 153. doi:10.1186/s13012-015-0345-7.
- AQ25** Ladner, D. 2015. *Strategy Testing: An Innovative Approach to Monitoring Highly Flexible Aid Programs (Case Study No 3), Working Politically in Practice*. Asia Foundation. 710
- Lang, T. 2011. "Adaptive Trial Design: Could We Use This Approach to Improve Clinical Trials in the Field of Global Health?." *American The Journal of tropical medicine and hygiene* 85: 967–970. doi:10.4269/ajtmh.2011.11-0151.
- Limato, R., R. Ahmed, A. Magdalena, S. Nasir, and F. Kotvojs. 2018. "Use of Most Significant Change (MSC) Technique to Evaluate Health Promotion Training of Maternal Community Health Workers in Cianjur District, Indonesia." *Evaluation and program planning* 66: 102–110. doi:10.1016/j.evalprogplan.2017.10.011. 715
- Lopez Bernal, J., S. Cummins, and A. Gasparrini. 2017. "Interrupted Time Series Regression for the Evaluation of Public Health Interventions: A Tutorial." *International journal of epidemiology* 46: 348–355. doi:10.1093/ije/dyw098.
- Lopez Bernal, J., S. Cummins, and A. Gasparrini. 2018. "The Use of Controls in Interrupted Time Series Studies of Public Health Interventions." *International journal of epidemiology*. doi:10.1093/ije/dyy135.
- Mahajan, R., and K. Gupta. 2010. "Adaptive Design Clinical Trials: Methodology, Challenges and Prospect." *Indian journal of pharmacology* 42: 201–207. doi:10.4103/0253-7613.68417. 720
- Manderson, L., and P. Aaby. 1992. "Can Rapid Anthropological Procedures Be Applied to Tropical Diseases?" *Health policy and planning* 7: 46–55. doi:10.1093/heapol/7.1.46.
- Manzano, A. 2016. "The Craft of Interviewing in Realist Evaluation." *Evaluation* 22: 342–360. doi:10.1177/1356389016638615. 725
- Mayne, J., 2008. "Contribution Analysis [WWW Document]." *Better Evaluation*. Accessed 26 October 2018. https://www.betterevaluation.org/en/plan/approach/contribution_analysis
- Moore, G. F., S. Audrey, M. Barker, L. Bond, C. Bonell, W. Hardeman, L. Moore, A. O'Cathain, T. Tinati, and D. Wight. 2015. "Process Evaluation of Complex Interventions: Medical Research Council Guidance." *BMJ (Clinical research ed.)* 350: h1258. doi:10.1136/bmj.h1258. 730
- O'Connell, T., and A. Sharkey. 2013. *Reaching Universal Health Coverage: Using a Modified Tanahashi Model Sub-Nationally to Attain Equitable and Effective Coverage*. New York: UNICEF.
- O'Donnell, M. 2016. *Adaptive Management: What It Means for Civil Society Organisations*. London: Bond.
- AQ26** ODI. 2009. *Strategy Development: Outcome Mapping*, In: *Tools for Knowledge and Learning: A Guide for Development and Humanitarian Organisations*. 735
- Optipedia, n.d. "A/B Testing [WWW Document]." *Optimizely*. Accessed 1 February 2018. <https://www.optimizely.com/optimization-glossary/ab-testing/>
- AQ27** Patton, M. Q. 2008. *Utilization-Focused Evaluation*. Sage publications.
- Patton, M. Q., 2013. "Utilization-Focused Evaluation (U-FE) Checklist [WWW Document]." Accessed 12 December 2017. https://wmich.edu/sites/default/files/attachments/u350/2014/UFE_checklist_2013.pdf 740

AQ28

Pawson, R. 2013. *The Science of Evaluation: A Realist Manifesto*. Sage.

AQ29

Pawson, R., and N. Tilley. 2004. *Realist Evaluation*.

Peerally, M. F., S. Carr, J. Waring, and M. Dixon-Woods. 2017. "The Problem with Root Cause Analysis." *BMJ quality & safety* 26: 417–422. doi:10.1136/bmjqs-2016-005511.

Peters, A., 2018. "At These Camps, Refugees Can Give Real-Time Customer Feedback [WWW Document]." *Fast Company*. Accessed 24 July 2018. <https://www.fastcompany.com/40575160/at-these-camps-refugees-can-give-real-time-customer-feedback> 745

Portela, M. C., P. J. Pronovost, T. Woodcock, P. Carter, and M. Dixon-Woods. 2015. "How to Study Improvement Interventions: A Brief Overview of Possible Study Types." *BMJ quality & safety* 24: 325–336. doi:10.1136/bmjqs-2014-003620.

Positive Deviance Initiative, 2017. "What Is Positive Deviance? [WWW Document]." *Positive Deviance Initiative*. URL 750
Accessed 20 March 2018. <https://positivedeviance.org/>

Research to Action, 2012. "Outcome Mapping: A Basic Introduction [WWW Document]." *Research to Action*. Accessed 7 December 2017. <http://www.researchtoaction.org/2012/01/outcome-mapping-a-basic-introduction/>

Rio, D., J. Hedges, S. Woodhead, and E. Rogers. 2015. *What Is the Bottleneck Analysis Approach for the Management of Severe Acute Malnutrition?*. UNICEF and Action Against Hunger. 755

AQ80

AQ81

Schünemann, H., Ed.. 2015. "Advances in Rapid Reviews." *Systematic reviews* 4.

Shiffman, S., A. A. Stone, and M. R. Hufford. 2008. "Ecological Momentary Assessment." *Annual review of clinical psychology* 4: 1–32. doi:10.1146/annurev.clinpsy.3.022806.091415.

Smutylo, T. 2005. "Outcome Mapping: A Method for Tracking Behavioural Changes in Development Programs (No.)" *ILAC Brief* 7. 760

AQ82

Srivastava, K., 2014. "The 'Adjacent Possible' of Big Data: What Evolution Teaches about Insights Generation [WWW Document]." *WIRED*. Accessed 21 January 2018. <https://www.wired.com/insights/2014/12/the-adjacent-possible-of-big-data/>

Stetler, C. B., M. W. Legro, C. M. Wallace, C. Bowman, M. Guihan, H. Hagedorn, B. Kimmel, N. D. Sharp, and J. L. Smith. 2006. "The Role of Formative Evaluation in Implementation Research and the QUERI Experience." *Journal of general internal medicine* 21: S1–S8. doi:10.1111/j.1525-1497.2006.00355.x. 765

Talcott, F., and V. Scholz. 2015. *Methodology Guide to Process Tracing for Christian Aid*. Oxford: International Non-Governmental Training and Research Centre.

Tanahashi, T. 1978. "Health Service Coverage and Its Evaluation." *Bulletin of the World Health Organization* 56: 295–303.

Theiss-Nyland, K., D. Koné, C. Karema, W. Ejersa, J. Webster, and J. Lines. 2017. "The Relative Roles of ANC and EPI in the Continuous Distribution of LLINs: A Qualitative Study in Four Countries." *Health policy and planning* 32: 467–475. doi:10.1093/heapol/czw158. 770

Thorlund, K., J. Haggstrom, J. J. Park, and E. J. Mills. 2018. "Key Design Considerations for Adaptive Clinical Trials: A Primer for Clinicians." *BMJ (Clinical research ed.)* 360: k698. doi:10.1136/bmj.k698.

Ton, G. 2012. "The Mixing of Methods: A Three-Step Process for Improving Rigour in Impact Evaluations." *Evaluation* 18: 5–25. doi:10.1177/1356389011431506. 775

Tricco, A. C., J. Antony, W. Zarin, L. Striffler, M. Ghassemi, J. Ivory, L. Perrier, B. Hutton, D. Moher, and S. E. Straus. 2015. "A Scoping Review of Rapid Review Methods." *BMC medicine* 13: 224. doi:10.1186/s12916-015-0465-6.

Tricco, A. C., E. Langlois, and S. E. Straus. 2017. *Rapid Reviews to Strengthen Health Policy and Systems: A Practical Guide*. Geneva: World Health Organization, Alliance for Health Policy and Systems Research. 780

Valters, C., C. Cummings, and H. Nixon. 2016. *Putting Learning at the Centre. Adaptive Development Programming in Practice*. Overseas Development Institute.

AQ83

Villar, S. S., J. Bowden, and J. Wason. 2017. "Response-Adaptive Designs for Binary Responses: How to Offer Patient Benefit while Being Robust to Time Trends?" *Pharmaceutical statistics*. 10.1002/pst.1845.

Vlassoff, C., and M. Tanner. 1992. "The Relevance of Rapid Assessment to Health Research and Interventions." *Health policy and planning* 7: 1–9. doi:10.1093/heapol/7.1.1. 785

Walji, A., and C. Vein, 2013. "Learning from Data-Driven Delivery [WWW Document]." World Bank. Accessed 11 October 2017. <http://blogs.worldbank.org/voices/learning-data-driven-delivery>

Webster, J., K. Kayentao, J. Bruce, S. I. Diawara, A. Abathina, A. A. Haiballa, O. K. Doumbo, and J. Hill. 2013. "Prevention of Malaria in Pregnancy with Intermittent Preventive Treatment and Insecticide Treated Nets in Mali: A Quantitative Health Systems Effectiveness Analysis." *PloS one* 8: e67520. doi:10.1371/journal.pone.0067520. 790

Wechsberg, W. M., J. W. Ndirangu, I. S. Speizer, W. A. Zule, W. Gumula, C. Peasant, F. A. Browne, and L. Dunlap. 2017. "An Implementation Science Protocol of the Women's Health CoOp in Healthcare Settings in Cape Town, South Africa: A Stepped-Wedge Design." *BMC women's health* 17. doi:10.1186/s12905-017-0433-8.

White, H., 2013. "Using the Causal Chain to Make Sense of the Numbers [WWW Document]." *International Initiative for Impact Evaluation*. Accessed 17 October 2018. <http://www.3ieimpact.org/en/announcements/2013/02/12/using-causal-chain-make-sense-numbers/> 795

White, H., and D. Phillips, 2012. "Addressing attribution of cause and effect in small n impact evaluations: towards an integrated framework (Working Paper 15)." International Initiative for Impact Evaluation.

Wilson-Grau, R., 2015. "Outcome Harvesting [WWW Document]." *Better Evaluation*. Accessed 4 January 2018 http://www.betterevaluation.org/en/plan/approach/outcome_harvesting 800

AQ84

- Wilson-Grau, R., and H. Britt. 2012. *Outome Harvesting Brief*. Ford Foundation.
- Wohlin, C., 2014. Guidelines for Snowballing in Systematic Literature Studies and a Replication in Software Engineering, in: Proceedings of the 18th International Conference on Evaluation and Assessment in Software Engineering, EASE '14. ACM, New York, pp. 38: 1–38:10. doi: [10.1145/2601248.2601268](https://doi.org/10.1145/2601248.2601268)
- Woolcock, M. 2009. "Toward a Plurality of Methods in Project Evaluation: A Contextualised Approach to Understanding Impact Trajectories and Efficacy." *Journal of development effectiveness* 1: 1–14. doi:[10.1080/19439340902727719](https://doi.org/10.1080/19439340902727719).

805

PROOF ONLY